

SPR Day 3

Readings

PRML 1.2.3 Bayesian probabilities

2.1 Binary rand. vars
+ Beta distrib.
+ Gamma functions

Goals

Review continuous rand. variables

Introduce Beta distribution

PDF formula

Graphs of PDF

What is a Gamma function?

Beta-Bernoulli joint model for coin toss

- posterior distribution
- MAP estimation
- predictive distribution & estimation
- evidence

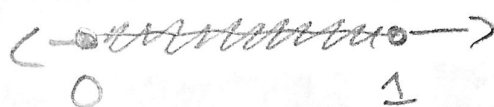
Continuous Random Variables

Examples

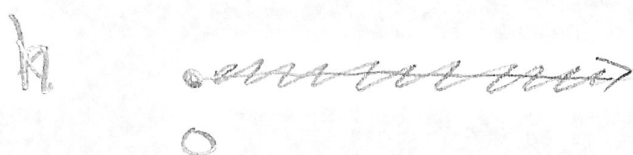
Probability coin ends up heads

μ sample space

$\mu \in [0, 1]$

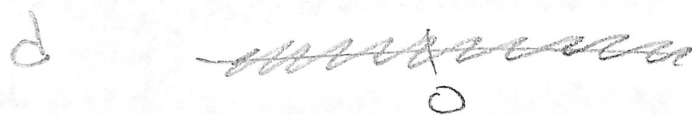


Height of avg. student



$h \in [0, \infty)$

Distance of random location from Medford



$d \in (-\infty, +\infty)$

Sample spaces are continuous intervals of real line

Can talk coherently about ^{the probability of} interval events

what's probability that height is above 75 cm

$$P(H \geq 75)$$

Need different language for specific events

should not say what is probability height equals 75
instead, think about limits

$$\text{pdf}(H=75) = P(H \in [75, 75+\delta]) \text{ as } \delta \rightarrow 0$$

↑
probability density function of random variable H

Two rules for a PDF function for continuous r.v. X

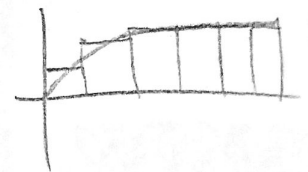
$$\text{pdf}(x) \geq 0 \quad \text{for all outcomes } x \in \mathcal{X}$$

$$\int_{x \in \mathcal{X}} \text{pdf}(x) dx = 1$$

Think of the Riemannian sum of this integral

$$\lim_{\Delta x \rightarrow 0} \sum_{n=1}^N \text{pdf}(x_n) \Delta x_n = 1$$

\uparrow \uparrow \uparrow
 density of volume mass
 each of of each bin
 bin each bin



$$\text{density} = \frac{\text{mass}}{\text{volume}}$$

Can density be larger than one? Yes, if mass of bin is small

Joint probabilities also obey these rules

$$X \in \mathbb{R}, Y \in \mathbb{R}$$

$$\text{pdf}(x, y) \geq 0 \quad \text{and} \quad \iint \text{pdf}(x, y) dx dy = 1$$

Sum Rule

$$\text{pdf}(x) = \int \text{pdf}(x, y) dy$$

Product Rule

$$\begin{aligned} \text{pdf}(x, y) &= \text{pdf}(x|y) \text{pdf}(y) \\ &= \text{pdf}(y|x) \text{pdf}(x) \end{aligned}$$

Beta distribution

Random variable μ

Sample space $\mu \in [0, 1]$

PDF function parameters $a > 0$
 $b > 0$

Beta pdf $(\mu/a, b)$

$$\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}$$

const. wrt μ

"interesting part" is a function of μ

we'll call this "normalizing constant"

$f(\mu/a, b)$

Gamma function: see notebook online

$$\Gamma(a) = \int_0^{\infty} u^{a-1} e^{-u} du$$

$$\Gamma: \mathbb{R} \rightarrow \mathbb{R}_+$$

scalar real \rightarrow scalar real

Remember:

$$\int_0^1 \text{BetaPDF}(\mu/a, b) d\mu = 1$$

$$\int_0^1 c(a, b) f(\mu/a, b) d\mu = 1$$

implies

$$\int_0^1 f(\mu/a, b) d\mu = \frac{1}{c(a, b)}$$

thus,

$$\int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

The Beta-Bernoulli Joint Model

We have defined a complete model (a joint distribution)

$$P(x_1, x_2, \dots, x_N, \mu) = \left[\prod_{n=1}^N \text{BernPMF}(x_n | \mu) \right] \cdot \text{BetaPDF}(\mu | a, b)$$

Bernoulli
"likelihood"

Beta
"prior"

Given this joint distribution, we can ask about related probability distributions that can be derived from this joint

"Posterior"

$$P(\mu | x_1, \dots, x_N) = \frac{P(x_1, \dots, x_N | \mu) P(\mu)}{P(x_1, \dots, x_N)}$$

"likelihood" "prior"

"evidence"

via Bayes rule

"Evidence"

$$P(x_1, \dots, x_N) = \int_{\mu=0}^1 P(x_1, \dots, x_N, \mu) d\mu$$

via sum rule

"Predictive Posterior"

$$P(x_N | x_1, \dots, x_{N-1}) = \int_{\mu=0}^1 P(x_N | \mu) P(\mu | x_1, \dots, x_{N-1}) d\mu$$

likelihood posterior

via sum rule & Bayes rule

Posterior of μ for Beta-Bern model

Using Bayes rule:

$$P(\mu | x_1, \dots, x_N) = \left[\prod_{n=1}^N \text{BernPMF}(x_n | \mu) \right] \cdot \text{BetaPDF}(\mu | a, b)$$

$$= \frac{1}{P(x_1, \dots, x_N)} \cdot \prod_{n=1}^N \mu^{x_n} (1-\mu)^{1-x_n} \cdot c(a, b) \cdot \mu^{a-1} (1-\mu)^{b-1}$$

Gather constant terms wrt μ

$$= \text{const} \cdot \mu^{(\sum_{n=1}^N x_n) + a - 1} \cdot (1-\mu)^{\sum_{n=1}^N (1-x_n) + b - 1}$$

Why? $(\mu^{x_1} \cdot \mu^{x_2} \cdot \mu^{x_3} \dots = \mu^{x_1 + x_2 + x_3 + \dots} = \mu^{\sum_{n=1}^N x_n})$

$$= \text{const} \cdot \mu^{\boxed{s(x) + a} - 1} \cdot (1-\mu)^{\boxed{r(x) + b} - 1}$$

$$= c(a', b') \cdot f(\mu, a', b') \quad \text{where } s(x) = \sum_{n=1}^N x_n$$

$$r(x) = N - s(x)$$

We have written the posterior PDF so it looks like " " " " "Beta" " " " " " " " "

constant wrt $\mu \cdot \mu^{a'-1} \cdot (1-\mu)^{b'-1}$

Since PDF must equal one when integrated, we have:

$$\frac{1}{c(a', b')} = \int_0^1 \mu^{a'-1} (1-\mu)^{b'-1} d\mu$$

$$c(a', b') = \frac{\Gamma(a' + b')}{\Gamma(a') \Gamma(b')}$$

Punchline:
 Posterior is Beta
 with parameters
 $a' = s(x) + a = \sum_{n=1}^N x_n + a$
 $b' = r(x) + b = \sum_{n=1}^N (1-x_n) + b$

Now, μ Towards MAP Estimator
and for Beta-Bernoulli x_1, \dots, x_N

$$P(\underbrace{x_1, x_2, \dots, x_N}_{\text{discrete}}, \underbrace{\mu}_{\text{continuous}}) = \text{Beta PDF}(\mu|a, b) \prod_{n=1}^N \text{Bern}(x_n|\mu)$$

We observe x_1, \dots, x_N , how to estimate μ ?
 how to predict next word x_{N+1} ?
 several ways:

- (1) find μ that maximizes $P(\mu|x_1, \dots, x_N)$
 called maximum a-posteriori (MAP)
- (2) take a more Bayesian approach (next time)

(1) MAP estimator

$$\begin{aligned} \mu^* &= \arg \max_{\mu \in [0,1]} P(\mu|x_1, \dots, x_N) \\ &= \arg \max_{\mu \in [0,1]} \frac{P(\mu) P(x_1, \dots, x_N|\mu)}{P(x_1, \dots, x_N)} \end{aligned}$$

const wrt μ
 $P(x_1, \dots, x_N) = \int_{\mu} P(\mu, x) d\mu$

$$= \arg \max_{\mu \in [0,1]} \log \text{Beta PDF}(\mu|a, b) + \sum_n \log \text{Bern PMF}(x_n|\mu)$$

$$= \arg \max_{\mu \in [0,1]} (a-1) \log \mu + (b-1) \log(1-\mu) + S(x) \log \mu + r(x) \log(1-\mu)$$

MAP estimator (cont'd)

$$\mu^* = \arg \max_{\mu \in [0,1]} \left(\overset{s'}{s(x)+a-1} \right) \log \mu + \left(\underset{r'}{r(x)+b-1} \right) \log (1-\mu)$$

We've solved this before! See derivation of ML estimator

$$\mu^* = \frac{s'}{s'+r'} \quad \text{when } \begin{matrix} s' \geq 0 \\ r' \geq 0 \end{matrix}$$

$$\mu^* = \frac{s(x)+a-1}{\underbrace{N+a+b-2}_w = s(x)+r(x)}$$

requires

$$a \geq 1$$

$$b \geq 1$$

otherwise
MAP does
not
exist

Interpretation

a acts like a "pseudocount" for heads

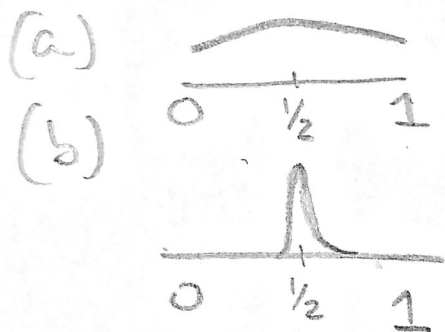
b acts like a "pseudocount" for tails

Problems w/ MAP estimator

- only works for $a > 1, b > 1$

What if I have $N=0, a=\frac{1}{2}, b=\frac{1}{2}$
can't make predictions.

- Should I really condense a whole distribution down to just one estimate?



should we act the same in these two cases?

Alternative

For predicting next coin (say X_{N+1} or X_*), consider:

ML $P(X_{N+1} | \mu^{ML}) = \frac{S(x)}{N}$

MAP $P(X_{N+1} | \mu^{MAP}) = \frac{S(x)+a-1}{N+a+b-2}$

Full Bayes

$P(X_{N+1} | X_1, \dots, X_N) = ? = \int P(X_{N+1} | \mu) P(\mu | X_1, \dots, X_N) d\mu$

condition on all data we observed

average over full distribution on μ

condition on reductive estimates of data

Posterior Predictive for next coin flip

Let's ask about probab. of heads:

$$P(X_{N+1}=1 | X_1, \dots, X_N) = \int \text{Bern}(x_{N+1}=1 | \mu) \text{Beta}(\mu | a', b') d\mu$$

$$= \int_0^1 \underbrace{\frac{\Gamma(a'+b')}{\Gamma(a')\Gamma(b')}}_{\text{const wrt } \mu} \mu^{a'-1} (1-\mu)^{b'-1} d\mu$$

$$= \frac{\Gamma(a'+b')}{\Gamma(a')\Gamma(b')} \int_0^1 \underbrace{\mu^{a'+1-1} (1-\mu)^{b'-1}}_{\text{combine like terms} = \mu^! \cdot \mu^{a'-1}} d\mu$$

const wrt μ
can be pulled out
of integral

combine like terms
 $= \mu^! \cdot \mu^{a'-1}$

recognise as unnormalized Beta PDF
use identity about

$$c(a', b') = \int f(\mu, a', b') d\mu$$

$$= \frac{\Gamma(a'+b')}{\Gamma(a')\Gamma(b')} \cdot \frac{\Gamma(a'+1)\Gamma(b')}{\Gamma(a'+1+b')}$$

cancel out
b' term

$$= \frac{\Gamma(a'+b')}{\Gamma(a')} \cdot \frac{a' \Gamma(a')}{(a'+b') \Gamma(a'+b')}$$

numerator use
identity
 $\Gamma(x+1) = x \Gamma(x)$
w/ $x=a$
denom. use same
w/ $x=a+b$

$$= \frac{a'}{a'+b'} = \boxed{\frac{S(x)+a}{N+a+b}}$$

This is the posterior predictive
works for all $a > 0, b > 0$